

THIS PROJECT: CROSS-SECTION OF STOCK RETURNS

Table: Literature surveys

Author(s)	Return predictive signals (RPS)
Subrahmanyam (2010)	50
Harvey, Liu and Zhu (2013)	185
McLean and Pontiff (2013)	82
Green, Hang and Zhang (2013)	330

[...] either US stock markets are pervasively inefficient, or there exist a much larger number of rationally priced sources of risk in equity returns than previously thought.

Green et al (2013)

TODAY: DISSECT MOMENTUM

- Usually modeled as stock return over past 12 months, skipping most recent month (M1-12)
- Novy-Marx (2012, JoFE) shows that an M7-12 strategy is more profitable than an M1-6 strategy
- Goyal and Wahal (2012) cannot find this effect in 36 out of 37 international markets
- What is going on? Which types of momentum can be learned from the data alone?

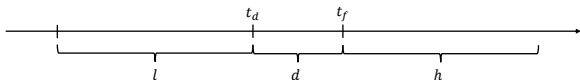
PERFORMANCE MEASURE

We measure performance by *precision*:

$$\text{Precision}(k) = \frac{\text{Correct predictions of class } k}{\text{All predictions of class } k}$$

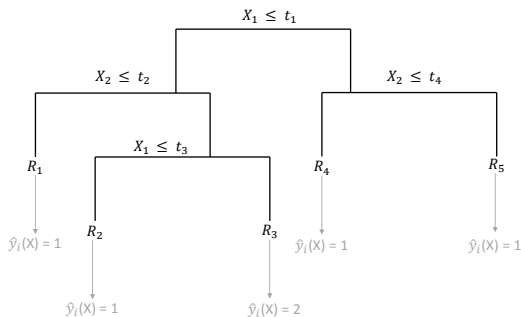
FEATURE CONSTRUCTION

- Construct momentum functions as $M(d, l)$
- $d \in \{0, \dots, 6\}$
- $l \in \{1, \dots, 18\}$
- $M(\cdot, \cdot)$ maps returns into quantiles \rightarrow 126 features/strategies
- More generally: $M(d, l, h)$. Today: $h = 1$.



DECISION-TREE ALGORITHM I

- Let y_i be a categorical outcome variable, and x_i a $p \times 1$ vector of candidate features.
- Basic idea with two features:



DECISION-TREE ALGORITHM II

- Finding optimal data partitioning is NP-complete → Use greedy procedure
- Share of class k in region R_m :

$$\hat{p}_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i = k)$$

- At each node, define *node impurity* as

$$H(\hat{p}) = - \sum_{k=1}^K \hat{p}_{mk} \log(\hat{p}_{mk})$$

- Find the feature j and the corresponding threshold that minimize $H(\hat{p})$ at each node

DECISION-TREE ALGORITHM III

- Automatic variable selection
- Robust to outliers
- Easily scalable to large N
- Prediction for new data point x is $\hat{T}(x)$.
 - Problem: Single tree is often a *weak learner*
 - Predictive accuracy low \rightarrow Random forests

ESTIMATION

- ① Decision-tree algorithm to deal with correlation of features
- ② **Random forest to reduce variance**
- ③ Hold-out sample to evaluate performance
- ④ Out-of-bag sampling to estimate *expected precision* and *feature importance* on test data

RANDOM FORESTS

Idea: Reduce variance of an estimate by averaging together many estimates

- Suppose B samples, identical but not independently distributed. $\hat{T}_b(x)$ is prediction of weak learner b at x . Let

$$\overline{\hat{T}_B(x)} = \frac{1}{B} \sum_{b=1}^B \hat{T}_b(x)$$

$$\text{Var}(\overline{\hat{T}_B(x)}) = \rho\sigma^2 + \frac{1-\rho}{B}\sigma^2 \rightarrow \rho\sigma^2, \text{ for } B \rightarrow \infty$$

- Random forests: Reduce ρ by training base learners on
 - Stochastic perturbation of the feature space
 - Random subset of the data

RANDOM FORESTS: ALGORITHM

Algorithm

For $b = 1 : B$

- ➊ Draw bootstrap sample Z of size N from training data
- ➋ Grow tree T_b that uses $m < p$ candidate features
- ➌ Collect $\{T_b\}_1^B$

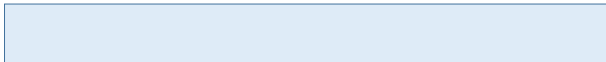
Class prediction is $\hat{C}_{RF}^B(x) = \text{majority vote}\{\hat{C}_b(x)\}_1^B$

RANDOM FORESTS: COMMENTS

- Bias of RF is the same as bias for any one tree → improvements in prediction are result of variance reduction alone
- Accuracy is often very high
- Random Forests lose the simple interpretability of a single tree

ESTIMATION

- ① Decision-tree algorithm to deal with correlation of features
- ② Random forest to reduce variance
- ③ **Hold-out sample to evaluate performance**
- ④ Out-of-bag sampling to estimate *expected precision* and *feature importance* on test data



ESTIMATION

- ① Decision-tree algorithm to deal with correlation of features
- ② Random forest to reduce variance
- ③ **Hold-out sample to evaluate performance**
- ④ Out-of-bag sampling to estimate *expected precision* and *feature importance* on test data

70%

30%

ESTIMATION

- ① Decision-tree algorithm to deal with correlation of features
- ② Random forest to reduce variance
- ③ Hold-out sample to evaluate performance
- ④ Out-of-bag sampling to estimate *expected precision* and *feature importance* on test data

RANDOM FORESTS: EVALUATION

Expected precision

- Bootstrap sampling with replacement leaves out $\approx 33\%$ of observations for each tree.
- Use those to estimate precision and average over all trees
- Compare to precision on hold-out sample



70%

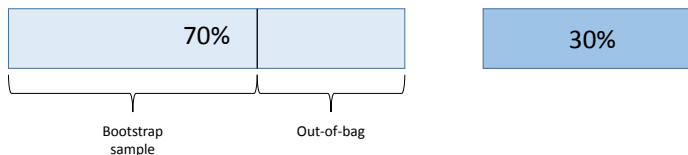


30%

RANDOM FORESTS: EVALUATION

Predictive precision

- Bootstrap sampling with replacement leaves out $\approx 33\%$ of observations for each tree.
- Use those to estimate the expected precision and average over all trees
- Compare to precision on hold-out sample



RANDOM FORESTS: EVALUATION

Feature importance

- 1 Compute prediction accuracy for OOB samples for T_b
- 2 Randomly permute values of variable j
- 3 Compute decrease in accuracy
- 4 Average over all trees

EMPIRICS

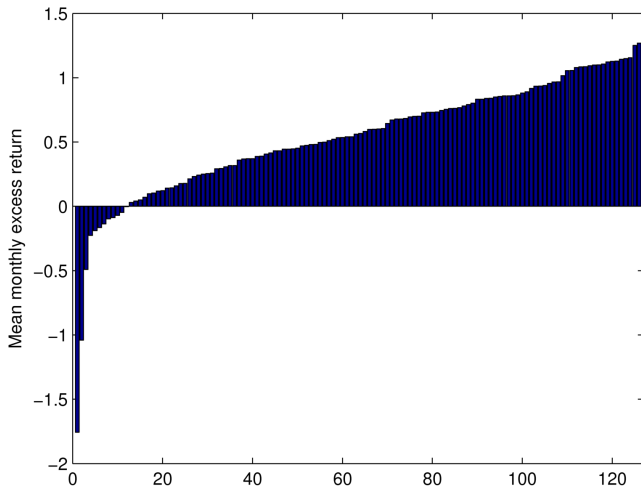
DATA

- CRSP monthly stock file, 1926-2012
- Exclude penny stocks
- Ordinary common shares only
- Truncate .5% and 99.5% monthly return percentile
- FF factors and UMD

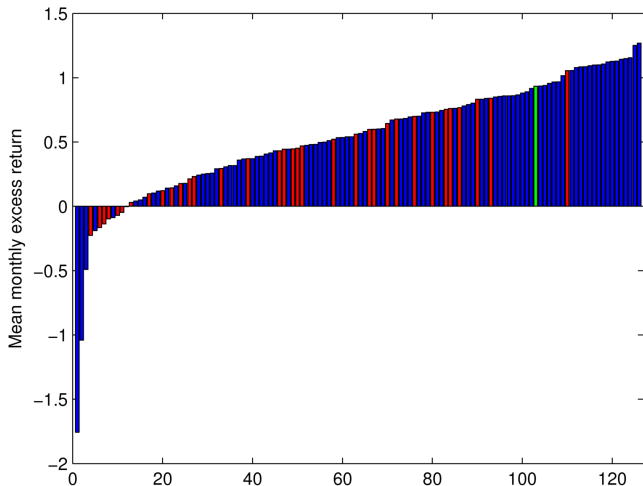
IMPLEMENTATION

- Construct 126 predictors $M(d, l)$
- Compute decile sorting for each predictor and one-month future return
- Tuning parameters (tree depth, number of randomly chosen predictors,...) can be cross-validated (not done yet)
- One RF run takes 6-7 hours with parallel processing

SINGLE FEATURE PERFORMANCE



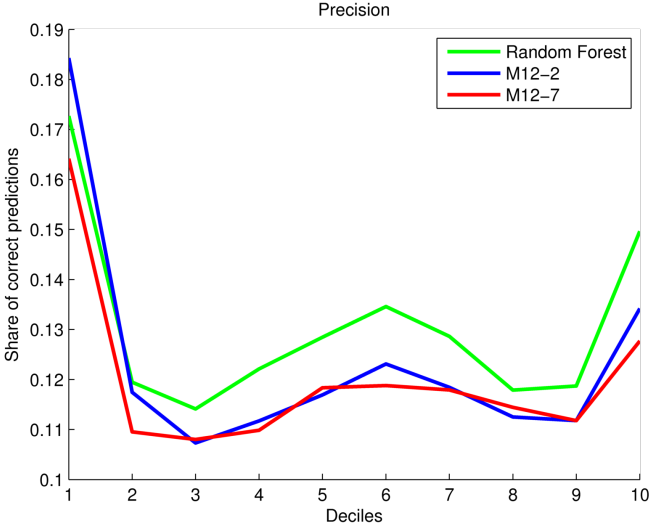
SINGLE FEATURE PERFORMANCE: MH ($d \in \{5, 6\}$)



IMPLEMENTATION

- Construct 126 predictors $M(d, l)$
- Compute decile sorting for each predictor and one-month future return
- Tuning parameters (tree depth, number of randomly chosen predictors,...) can be cross-validated (not done yet)
- One RF run takes 6-7 hours with parallel processing

RANDOM FOREST: PRECISION



RANDOM FOREST: FEATURE IMPORTANCE

d (lag)	l (length)	Relative Importance
0	1	1
0	2	.660
2	1	.655
1	1	.595
0	3	.555
2	3	.555
1	3	.542
1	2	.541
4	2	.527
3	1	.502

RANDOM FOREST: FEATURE IMPORTANCE

d (lag)	l (length)	Relative Importance
0	1	1
0	2	.660
2	1	.655
1	1	.595
0	3	.555
2	3	.555
1	3	.542
1	2	.541
4	2	.527
3	1	.502

RANDOM FOREST: FEATURE IMPORTANCE

Most important		Least important	
d (lag)	l (length)	d (lag)	l (length)
0	1	5	12
0	2	2	12
2	1	3	11
1	1	2	17
0	3	0	17
2	3	0	14
1	3	2	15
1	2	3	16
4	2	2	13
3	1	3	12

RANDOM FOREST: FEATURE IMPORTANCE

Most important		Least important	
d (lag)	l (length)	d (lag)	l (length)
0	1	5	12
0	2	2	12
2	1	3	11
1	1	2	17
0	3	0	17
2	3	0	14
1	3	2	15
1	2	3	16
4	2	2	13
3	1	3	12

Brief look at simple strategy return:

- Buy highest predicted quantile
- Sell lowest predicted quantile

RANDOM FOREST: FACTOR LOADINGS

Strategy: Long highest predicted quantile, short lowest predicted quantile

	(1)	(2)	(3)	(4)
Intercept	.97 (5.60)	1.14 (5.56)	1.26 (5.86)	.76 (3.67)
MKT		-.18 (-1.80)	-.10 (-1.21)	.01 (.18)
SMB			-.18 (-1.58)	-.16 (-1.54)
HML			-.36 (-3.70)	-.14 (-1.55)
UMD				.47 (8.02)

RANDOM FOREST: FACTOR LOADINGS

Strategy: Long highest predicted quantile, short lowest predicted quantile

	(1)	(2)	(3)	(4)
Intercept	.97 (5.60)	1.14 (5.56)	1.26 (5.86)	.76 (3.67)
MKT		-.18 (-1.80)	-.10 (-1.21)	.01 (.18)
SMB			-.18 (-1.58)	-.16 (-1.54)
HML			-.36 (-3.70)	-.14 (-1.55)
UMD				.47 (8.02)

EXCLUDING FEATURES

- What is the relevance of whole classes of features?
- Idea: Does exclusion of one group result in performance deterioration in the prediction task?
- Precision might even be higher if focus on more relevant features
- Define short-horizon features as $d \leq 1$, medium-horizon as $d \geq 5$

MOST IMPORTANT FEATURES

Rank	$\{X : d \leq 4\}$		$\{X : d \geq 2\}$	
	d (lag)	l (length)	d (lag)	l (length)
1	0	1	2	1
2	0	2	2	2
3	1	1	5	1
4	1	2	2	3
5	0	3	6	1

MOST IMPORTANT FEATURES

Rank	$\{X : d \leq 4\}$		$\{X : d \geq 2\}$	
	d (lag)	l (length)	d (lag)	l (length)
1	0	1	2	1
2	0	2	2	2
3	1	1	5	1
4	1	2	2	3
5	0	3	6	1

FACTOR LOADINGS

	(1)	$\{X : d \leq 4\}$		(4)	(5)	$\{X : d \geq 2\}$		(8)
		(2)	(3)			(6)	(7)	
Intercept	1.34	1.46	1.59	1.08	.82	.97	1.08	.48
	(6.42)	(6.49)	(6.94)	(5.04)	(4.55)	(5.00)	(5.18)	(2.38)
MKT		-.15	-.10	-.02		-.16	-.09	.02
		-(1.92)	-(1.55)	-(.34)		-(2.44)	-(1.35)	(.28)
SMB			-.09	-.05			-.21	-.20
			-(1.03)	-(.67)			-(2.30)	-(2.73)
HML			-.38	-.19			-.28	-.05
			-(3.71)	-(2.11)			-(2.71)	-(.46)
UMD				.47				.57
				(8.17)				(10.59)

FACTOR LOADINGS

	$\{X : d \leq 4\}$				$\{X : d \geq 2\}$			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Intercept	1.34 (6.42)	1.46 (6.49)	1.59 (6.94)	1.08 (5.04)	.82 (4.55)	.97 (5.00)	1.08 (5.18)	.48 (2.38)
MKT		-.15 (-1.92)	-.10 (-1.55)	-.02 (-.34)		-.16 (-2.44)	-.09 (-1.35)	.02 (.28)
SMB			-.09 (-1.03)	-.05 (-.67)			-.21 (-2.30)	-.20 (-2.73)
HML			-.38 (-3.71)	-.19 (-2.11)			-.28 (-2.71)	-.05 (-.46)
UMD				.47 (8.17)				.57 (10.59)

SUMMARY: EXCLUDING FEATURES

- Eyeballing suggests that precision suffered by excluding SH features, and did not change excluding MH features
- Return regressions point in the same direction
- Need to test this formally
- Additional test: Predict SH features with MH features (and vice versa) and use fitted values to predict future returns (flavor of Granger causality test)

